

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : **2003-271616**

(43)Date of publication of application : **26.09.2003**

(51)Int.Cl.

**G06F 17/30**

(21)Application number : **2002-068858**

(71)Applicant : **RICOH CO LTD**

(22)Date of filing : **13.03.2002**

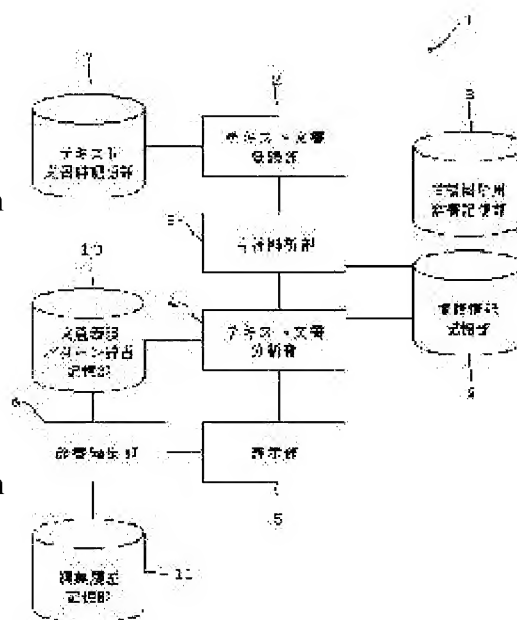
(72)Inventor : **SATO NAOKO**

## (54) DOCUMENT CLASSIFICATION DEVICE, DOCUMENT CLASSIFICATION METHOD AND RECORDING MEDIUM

(57)Abstract:

**PROBLEM TO BE SOLVED:** To provide a document classification device, document classification method and recording medium for classifying a text document on the basis of a desired meaning.

**SOLUTION:** In this document classification device 1, a language analysis part 3 retrieves a text document including an inputted retrieval request from a plurality of text documents collected and accumulated in a text document group storage part 2, and classifies the plurality of text documents by categories on the basis of the text document retrieval result. A text document analysis part 4 classifies, on the basis of the category classification result, the text documents by meanings in reference to the meaning expression pattern dictionary of a meaning expression pattern dictionary storage part 10 in which characteristic meaning expression patterns for expressing meanings are registered, and a dictionary editing part 6 optionally edits the meaning expression pattern dictionary according to a user's operation. Accordingly, the classification reference can be optionally changed to perform the document classification according to the meaning desired by the user, and the convenience can be improved.



## LEGAL STATUS

[Date of request for examination] 01.03.2005

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号  
特開2003-271616  
(P2003-271616A)

(43)公開日 平成15年9月26日(2003.9.26)

(51)Int.Cl. <sup>7</sup>	識別記号	F I	テーマコード <sup>*</sup> (参考)
G 0 6 F 17/30	2 1 0	G 0 6 F 17/30	2 1 0 D 5 B 0 7 5
	1 7 0		1 7 0 A
	2 2 0		2 2 0 Z

審査請求 未請求 請求項の数7 O L (全 11 頁)

(21)出願番号 特願2002-68858(P2002-68858)

(22)出願日 平成14年3月13日(2002.3.13)

(71)出願人 000006747

株式会社リコー

東京都大田区中馬込1丁目3番6号

(72)発明者 佐藤 奈穂子

東京都大田区中馬込1丁目3番6号 株式  
会社リコー内

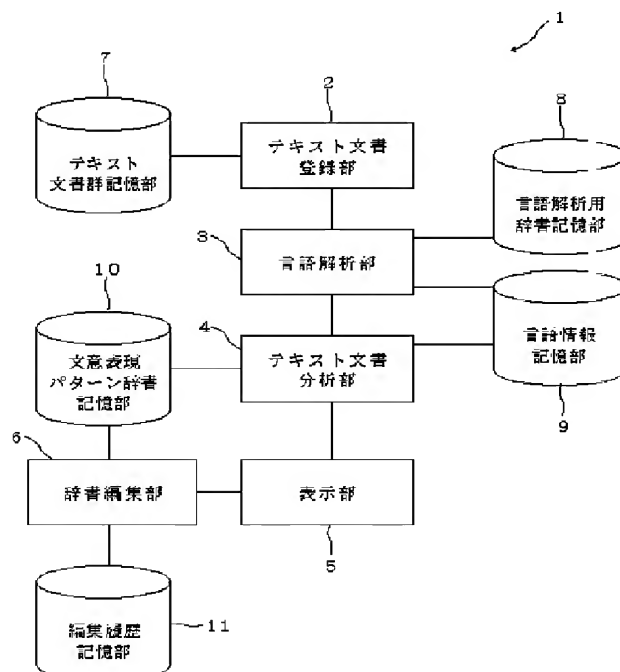
Fターム(参考) 5B075 ND03 ND34 NK06 NK32 NK39  
NR02 NR12 NS10 UU06

(54)【発明の名称】 文書分類装置、文書分類方法及び記録媒体

(57)【要約】

【課題】本発明はテキスト文書を所望の文意による文書分類する文書分類装置、文書分類方法及び記録媒体を提供する。

【解決手段】文書分類装置1は、言語解析部3で、入力された検索要求を含むテキスト文書を、テキスト文書群記憶部2に収集・蓄積された複数のテキスト文書から検索し、言語解析部3で、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、テキスト文書分析部4が、当該カテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書記憶部10の文意表現パターン辞書を参照してテキスト文書を文意別に分類し、文意表現パターン辞書を、ユーザの操作に応じて、辞書編集部6で、任意に編集する。したがって、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行うことができ、利用性を向上させることができる。



## 【特許請求の範囲】

【請求項1】入力された検索要求を含むテキスト文書を文書蓄積手段に収集・蓄積された複数のテキスト文書から検索するテキスト文書検索手段と、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成手段と、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を記憶する文意表現パターン辞書記憶手段と、前記カテゴリ生成手段のカテゴリ分け結果に基づいて前記文意表現パターン辞書を参照して前記テキスト文書を文意別に分析・分類する文意分類手段と、を備えた文書分類装置であって、前記文意表現パターン辞書記憶手段の文意表現パターン辞書を任意に編集する文意表現パターン辞書編集手段を備えていることを特徴とする文書分類装置。

【請求項2】前記文意表現パターン辞書は、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集手段は、前記文意表現パターン辞書に、前記意図タグと前記文意表現パターンを新たに追加登録することを特徴とする請求項1記載の文書分類装置。

【請求項3】前記文意表現パターン辞書は、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集手段は、前記文意表現パターン辞書に登録されている前記意図タグと前記文意表現パターンを個々に無効化することを特徴とする請求項1または請求項2記載の文書分類装置。

【請求項4】文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類する文書分類方法において、入力された検索要求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行うことを特徴とする文書分類方法。

【請求項5】前記文書分類方法は、前記文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集処理ステップが、前記文意表現パターン辞書に、前記意図タグと前記文意表現パターンを新たに追加登録することを特徴とする請求項4記載の文書分類方法。

【請求項6】前記文書分類方法は、前記文意表現パター

ン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集処理ステップが、前記文意表現パターン辞書に登録されている前記意図タグと前記文意表現パターンを個々に無効化することを特徴とする請求項1または請求項2記載の文書分類方法。

【請求項7】文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類する文書分類方法のプログラム及びデータを記憶する記録媒体であって、入力された検索要求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行う文書分類方法のプログラム及びデータを記録することを特徴とする記録媒体。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】本発明は、文書分類装置、文書分類方法及び記録媒体に関し、詳細には、テキスト文書をカテゴリ分けし、文意表現パターン辞書を参照して文意別に分析・分類するとともに、文意表現パターン辞書を任意に編集し、所望の文意による文書分類を行う文書分類装置、文書分類方法及び当該文書分類方法のプログラムとデータを記録した記録媒体に関する。

## 【0002】

【従来の技術】近時、情報の電子化が進み、従来紙文書で保管されていた文書も電子化されるようになってきている。このような文書の電子化に伴って、大量の電子化文書が流通し、収集・蓄積された電子化文書をいかに管理して簡便に再利用するかが重量な問題となってきた。

【0003】そして、従来、このような電子文書の簡便な利用や管理のためにさまざまな文書処理技術が提案されており、この文書処理技術の一例として、ある目的で収集された文書群の自動分類が挙げられる。この文書群の自動分類は、大量の電子化文書群から類似した文書群を自動分類する技術であり、一般的には、各文書に含まれている重要語句の類似性、出現頻度、出現場所等の共通点に基づいて、関連性の高い文書をグルーピングする仕組みになっている。この文書群の自動分類で利用されている重要語句としては、従来、文書におけるキーワードが用いられており、主に文書中に頻出する名詞、動詞等の品詞を限定して抽出した単語である。

【0004】また、従来、処理対象となる文書を入力する入力部と、入力した文書中の文章に対して形態素解析を行う形態素解析部と、前記形態素解析部から出力された形態素列の部分列を、重み付きで特定表現候補として取得する特定表現候補取得部と、予めいくつかの特定表現を格納した特定表現辞書と、与えられた形態素列の前記特定表現辞書中の表現に対するマッチ度を表す実数を、当該形態素列の前記特定表現辞書に対する検索結果として取得する特定表現辞書検索部と、前記特定表現候補に対して、前記候補に付与された重みと、前記特定表現検索部による前記候補の前記特定表現辞書に対する検索結果とを変数として判別スコアを計算し、前記判別スコアが予め設定した一定の値を下回る候補を除外する判別分析実行部と、前記特定表現候補のうち、前記判別分析実行部によって除外されなかった形態素の文字列を特定表現として出力する出力部とを備えた文書処理装置が提案されている（特開2001-75959号公報参照）。

【0005】すなわち、この従来技術は、人名や企業名等に特有の特定表現を、辞書を用いて抽出し、それを重要語句として、単語単体よりも、複数の単語の並びや単語の出現パターンなど含有情報が多く、より限定がかった表現単位での類似性判断を行なうことで、分類の精度の向上を図っている。

【0006】

【発明が解決しようとする課題】しかしながら、上記公報記載の従来技術は、人名や企業名等の特定表現の情報抽出を目的としたものであり、特定表現を利用した分類処理については、言及されておらず、また、特定表現パターンの内容として、書き手の文意を表わす表現パターンについては、上記公報の実施例に記載されていない。

【0007】一方、アンケートの自由記述部分などを分類、分析する場合、設問に関する話題は予め分かっており、書き手の意図にこそ分析のポイントがあると考えられ、書き手の文意を表わす表現パターンが要求されることは必至である。また、その際に、意図表現パターンをユーザが目的に応じて解釈し、分類（分析）基準を任意に変更できるようなくみが要望されている。

【0008】そこで、請求項1記載の発明は、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類するに際して、入力された検索要求を含むテキスト文書を、テキスト文書検索手段で、文書蓄積手段に収集・蓄積された複数のテキスト文書から検索し、カテゴリ生成手段で、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、カテゴリ生成手段のカテゴリ分け結果に基づいて、文意分類手段が、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照してテキスト文書を文意別に分析・分類し、当該文意表現パターン辞書を、文意表現パターン辞書編集手段で、任意に編集することにより、

分類基準を任意に変更して、ユーザの所望する文意により文書分類を行い、利用性の良好な文書分類装置を提供することを目的としている。

【0009】請求項2記載の発明は、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集手段が、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録することにより、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行い、利用性の良好な文書分類装置を提供することを目的としている。

【0010】請求項3記載の発明は、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集手段が、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化することにより、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにし、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性の良好な文書分類装置を提供することを目的としている。

【0011】請求項4記載の発明は、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類するに際して、テキスト文書検索処理ステップで、入力された検索要求を含むテキスト文書を文書蓄積手段のテキスト文書から検索し、カテゴリ生成処理ステップで、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、文意分類処理ステップで、カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書を任意に編集することにより、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行い、利用性の良好な文書分類方法を提供することを目的としている。

【0012】請求項5記載の発明は、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録することにより、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行い、利用性の良好な文書分類方法を提供することを目的としている。

【0013】請求項6記載の発明は、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文

書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化することにより、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにし、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性の良好な文書分類方法を提供することを目的としている。

【0014】請求項7記載の発明は、文書蓄積手段に収集・蓄積された複数のテキスト文書进行分析・分類する文書分類方法のプログラム及びデータを記憶する記録媒体に、当該プログラム及びデータとして、入力された検索要求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行う文書分類方法のプログラム及びデータを記録することにより、当該プログラム及びデータをコンピュータ等の情報処理装置に導入することで、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行い、利用性を向上させることのできる文書分類装置及び文書分類方法を実現する記録媒体を提供することを目的としている。

【0015】

【課題を解決するための手段】請求項1記載の発明の文書分類装置は、入力された検索要求を含むテキスト文書を文書蓄積手段に収集・蓄積された複数のテキスト文書から検索するテキスト文書検索手段と、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成手段と、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を記憶する文意表現パターン辞書記憶手段と、前記カテゴリ生成手段のカテゴリ分け結果に基づいて前記文意表現パターン辞書を参照して前記テキスト文書を文意別に分析・分類する文意分類手段と、を備えた文書分類装置であって、前記文意表現パターン辞書記憶手段の文意表現パターン辞書を任意に編集する文意表現パターン辞書編集手段を備えていることにより、上記目的を達成している。

【0016】上記構成によれば、文書蓄積手段に収集・蓄積された複数のテキスト文書进行分析・分類するに際して、入力された検索要求を含むテキスト文書を、テキス

ト文書検索手段で、文書蓄積手段に収集・蓄積された複数のテキスト文書から検索し、カテゴリ生成手段で、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、カテゴリ生成手段のカテゴリ分け結果に基づいて、文意分類手段が、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照してテキスト文書を文意別に分析・分類し、当該文意表現パターン辞書を、文意表現パターン辞書編集手段で、任意に編集するので、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行うことができ、利用性を向上させることができる。

【0017】この場合、例えば、請求項2に記載するように、前記文意表現パターン辞書は、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集手段は、前記文意表現パターン辞書に、前記意図タグと前記文意表現パターンを新たに追加登録するものであってもよい。

【0018】上記構成によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集手段が、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録するので、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行うことができ、利用性を向上させることができる。

【0019】また、例えば、請求項3に記載するように、前記文意表現パターン辞書は、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集手段は、前記文意表現パターン辞書に登録されている前記意図タグと前記文意表現パターンを個々に無効化するものであってもよい。

【0020】上記構成によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集手段が、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化するので、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにすることができ、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性を向上させることができる。

【0021】請求項4記載の発明の文書分類方法は、文書蓄積手段に収集・蓄積された複数のテキスト文書进行分析・分類する文書分類方法において、入力された検索要

求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行うことにより、上記目的を達成している。

【0022】上記構成によれば、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類するに際して、テキスト文書検索処理ステップで、入力された検索要求を含むテキスト文書を文書蓄積手段のテキスト文書から検索し、カテゴリ生成処理ステップで、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、文意分類処理ステップで、カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書を任意に編集するので、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行うことができ、利用性を向上させることができる。

【0023】この場合、例えば、請求項5に記載するように、前記文書分類方法は、前記文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集処理ステップが、前記文意表現パターン辞書に、前記意図タグと前記文意表現パターンを新たに追加登録してもよい。

【0024】上記構成によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録するので、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行うことができ、利用性を向上させることができる。

【0025】また、例えば、請求項6に記載するように、前記文書分類方法は、前記文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとが対として登録されており、前記文意表現パターン辞書編集処理ステップが、前記文意表現パターン辞書に登録されている前記意図タグと前記文意表現パターンを個々に

無効化してもよい。

【0026】上記構成によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化するので、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにすることができ、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性を向上させることができる。

【0027】請求項7記載の発明の記録媒体は、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類する文書分類方法のプログラム及びデータを記憶する記録媒体であって、入力された検索要求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行う文書分類方法のプログラム及びデータを記録することにより、上記目的を達成している。

【0028】上記構成によれば、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類する文書分類方法のプログラム及びデータを記憶する記録媒体に、当該プログラム及びデータとして、入力された検索要求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行う文書分類方法のプログラム及びデータを記録しているので、当該プログラム及びデータをコンピュータ等の情報処理装置に導入することで、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行い、利用性を向上させることのできる文書分類装置及び文書分類方法を実現することができる。

【0029】



【発明の実施の形態】以下、本発明の好適な実施の形態を添付図面に基づいて詳細に説明する。なお、以下に述べる実施の形態は、本発明の好適な実施の形態であるから、技術的に好ましい種々の限定が付されているが、本発明の範囲は、以下の説明において特に本発明を限定する旨の記載がない限り、これらの態様に限られるものではない。

【0030】図1～図4は、本発明の文書分類装置、文書分類方法及び記録媒体の一実施の形態を示す図であり、図1は、本発明の文書分類装置、文書分類方法及び記録媒体の一実施の形態を適用した文書分類装置1のブロック構成図である。

【0031】図1において、文書分類装置1は、テキスト文書登録部2、言語解析部3、テキスト文書分析部4、表示部5、辞書編集部6、テキスト文書群記憶部7、言語解析用辞書記憶部8、言語情報記憶部9、文意表現パターン辞書記憶部10及び編集履歴記憶部11等を備えている。文書分類装置1は、文書分類処理プログラム及び必要なデータを記録するCD-ROM (Compact Disc Read Only Memory) 等の記録媒体を、例えば、コンピュータ等に読み取らせて導入することで、構築される。

【0032】テキスト文書群記憶部(文書蓄積手段)7は、収集されたテキスト文書のテキスト文書群が登録され、登録されたテキスト文書が文書分類・分析の対象となる。

【0033】テキスト文書登録部2は、収集されたテキスト文書をテキスト文書群記憶部7に登録して蓄積させ、登録したテキスト文書群の管理を行う。

【0034】言語解析用辞書記憶部8は、言語解析部3による言語解析に必要な各種言語解析情報を記憶する。

【0035】言語情報記憶部9は、言語解析部3によるテキスト文書の解析処理によって得られる言語的属性を解析単位毎に記憶する。

【0036】言語解析部(テキスト文書検索手段、カテゴリ生成手段)3は、入力された検索要求を含むテキスト文書をテキスト文書群記憶部7に収集・蓄積された複数のテキスト文書から検索するテキスト文書検索処理と、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理を行う。具体的には、言語解析部3は、言語解析用辞書記憶部8の記憶する言語解析用辞書に基づいて、テキスト文書登録部2によりテキスト文書群記憶部7に登録された各テキスト文書の形態素を解析する形態素解析処理、形態素解析処理の解析結果と文節生成規則に基づいて文節を生成する文節生成処理、文節生成処理で生成した文節が解析対象のテキスト文書のどの文節に係っているかを予め設定されている係り受け解析規則に基づいて特定する係り受け解析処理等の各ステップ処理を実行し、これらの形態素解析処理、文節生成処理、係り受け解析処理によっ

て得られる言語的属性を解析単位に言語情報記憶部9に保持させる言語情報登録処理を行う。

【0037】文意表現パターン辞書記憶部(文意表現パターン辞書記憶手段)10は、文意を表現する特徴的な文意表現パターン辞書が登録され、蓄積する。この文意表現パターン辞書記憶部10に登録される文意表現パターン辞書としては、例えば、図2に示すようなものであり、文書の書き手の意図を示す文意タグの種類(例えば、依頼、要望、否定、疑問、推量等)毎に、文意表現パターンがそれぞれ対として複数登録されている。例えば、依頼の文意タグに対しては、「用言+て/欲しい」、「用言+て/ください」、「サ変名詞+ください」等の文書表現パターン、要望の文意タグに対しては、「用言+助動詞たい」、「用言+て/いただき・助動詞たい」、「サ変名詞+いただき・助動詞たい」等の文書表現パターン、否定の文意タグに対しては、「用言+助動詞ない」、「サ変名詞+助動詞+助動詞ない」等の文書表現パターン、疑問の文意タグに対しては、「用言+終助詞か」、「文末=記号?」等の文書表現パターン、推量の文意タグに対しては、「用言+助動詞だろう」、「用言+副助詞か+副助詞も+しれない」、「用言+助動詞そう」、「助言+助動詞できる+助動詞そう」等の文書表現パターンが対として登録されている。

【0038】辞書編集部(文意表現パターン辞書編集手段)6は、文意表現パターン辞書記憶部10に登録されている文意表現パターン辞書を、ユーザの操作に応じて、任意に編集し、例えば、文意タグと文意表現パターンを文意表現パターン辞書に新たに追加登録し、また、文意表現パターン辞書に記載されている文書の書き手の意図を示す文意タグと、それを特定するための特徴的な言語表現を個々に無効化する。

【0039】編集履歴記憶部11は、辞書編集部6による文意表現パターン辞書記憶部10の文意表現パターン辞書の編集履歴を記憶する。

【0040】テキスト文書分析部(文書分類手段)4は、言語解析部3による解析結果の各解析単位または各単位における言語的属性に基づいて、テキスト文書から文意表現パターン辞書記憶部10の文意表現パターン辞書を検索可能な単位を生成し、この文意表現パターン辞書を検索可能な単位として、例えば、各テキスト文書に対して、文節を生成する文意別文節処理を行う。また、テキスト文書分析部4は、日本語では、文末に意図表現が表出するという特性を利用するために、文末文節を抽出する。

【0041】すなわち、テキスト文書分析部4は、言語解析部3の解析結果である各解析単位または各単位における言語属性を、文書表現パターン辞書記憶部8の文書表現パターンと文意タグを検索可能な形式に変換する文意表現変換処理、当該文意表現変換処理した言語属性に基づいて文書表現パターン辞書記憶部8の辞書引きを行



う辞書引き処理、当該辞書引き処理で辞書引きした文書表現パターンに合致した文書表現パターンを文意タグに変換する文意タグ変換処理及び文意タグ変換処理で変換した文意タグを用いてテキスト文書の文意別カテゴリの分類を行う文意別カテゴリ分類処理を行う。

【0042】表示部5は、液晶ディスプレイやCRT（陰極線管：Cathode Ray Tube）等が用いられ、テキスト文書分析部4の分析したり、分類した結果を表示し、また、辞書編集部6の編集結果等を表示する。

【0043】次に、本実施の形態の作用を説明する。文書分類装置1は、文書分類処理プログラム及び必要なデータを記録するCD-ROM等の記録媒体を、例えば、コンピュータ等の情報処理装置に読み取らせて導入することで、構築され、電子化されたテキスト文書群を言語解析して、書き手の意図を表現する表現パターンをテキスト文書中から重要語句として抽出し、文書を分類するところにその特徴がある。

【0044】すなわち、文書分類装置1は、分析対象のテキスト文書群が入力されると、当該テキスト文書群をテキスト文書登録部2が当該テキスト文書群をテキスト文書群記憶部5に登録する。

【0045】いま、例えば、あるマリンスポーツについて意見を収集・蓄積した、下記のようなテキストデータがあり、集めた意見を文意別に分類・分析するものとして、以下説明する。

～ 収集テキスト ～

1. お金がかかりそう。
2. サーフィンをやれる環境をもっと良くしてほしい！
3. もっと盛んになってほしい。
4. 全てのコトが忘れられて、すごく楽しそう。
5. 安く手軽にできるならやってみたい。
6. リフレッシュできそう。
7. もっと海岸でのマナーを大切に指導して欲しい。
8. もっと安くして欲しい。

1. 用言+助動詞そう
2. 用言+て／ほしい
3. 用言+て／ほしい
4. 用言+助動詞そう
5. 用言+助動詞たい
6. 用言+助動詞できる+助動詞そう
7. 用言+て／欲しい
8. 用言+て／欲しい
9. 用言+助動詞そう
10. 用言+助動詞そう
11. 用言+助動詞そう

この辞書引きの結果、文意（推量）は、{1, 4, 6, 9, 10, 11}、文意（依頼）は、{2, 3, 7, 8}、文意（要望）は、{5}という文書の文意別カテゴリ分類を行うことができる。

【0050】そして、文書分類装置1は、上記文書分類

9. ととても楽しそうだが面倒くさそう。

10. お金がいっぱいかかるけど楽しそう。

11. 一回やったらハマリそう。

【0046】まず、最初に、これらのテキスト文書に対して、言語解析部3で言語解析を行い、テキストの構成単語の品詞等の属性情報を取得する。この言語解析は、既存のさまざまな手法で実現することができる。

【0047】次に、テキスト文書分析部4が、言語解析部3の解析結果に基づいて、文意表現パターン辞書記憶部10の文意表現パターン辞書を検索可能な単位を生成、例えば、各テキストに対して、文節を単位として、生成する。この文節生成技術は、既存の言語処理技術で実現することができる。

【0048】次に、テキスト文書分析部4は、日本語では文末に意図表現が表出するという特性を利用して、各テキスト文書について、以下のように、文末文節を抽出する。

1. かかり・そう・。
2. 良く・し・て・ほしい・！
3. なっ・て・ほしい・。
4. 楽し・そう・。
5. やっ・て・み・たい・。
6. リフレッシュ・でき・そう・。
7. 指導・し・て・欲しい・。
8. し・て・欲しい・。
9. 面倒くさ・そう・。
10. 楽し・そう・。
11. ハマリ・そう・。

【0049】テキスト文書分析部4は、この単位で、品詞等の属性を用いて、さらに正規表現に変換して記憶し、これらに対して、図2に示した文意表現パターン辞書を辞書引きする。この辞書引きにより、上記例では、以下のような辞書引き結果を得ることができる。

- 文意（推量）
- 文意（依頼）
- 文意（依頼）
- 文意（推量）
- 文意（要望）
- 文意（推量）
- 文意（依頼）
- 文意（依頼）
- 文意（推量）
- 文意（推量）
- 文意（推量）

結果を得たユーザが、上記テキスト文書の例の場合に、文意（依頼）は文意（要望）に併せたほうが適切であると考えられる場合、辞書編集機能を利用して、文意（依頼）の文意表現パターンの無効化や文意（要望）の文意表現パターンの新規追加を行なうことができる。

【0051】すなわち、ユーザが文書分類装置1の図示しない操作部で、ボタンやコマンド入力等を行って、辞書編集の指示を出すと、文書分類装置1は、表示部5の表示画面に、例えば、図3に示すような辞書編集ウィンドウを表示させる。文書分類装置1は、この辞書編集ウィンドウとして、図3に示すように、文意表現パターン辞書記憶部10の文意表現パターン辞書に登録されている文意と文意表現パターン対を表示し、また、各文意表現パターン対に、無効化欄が設けられている。

【0052】ユーザは、この文意表現パターン対の無効化欄に、図3に示すように、チェックを入れることで、当該文意の対として登録されている文意表現パターンを個別に無効化することができ、文書分類装置1は、無効化欄にチェックの入れられた文意の対として登録されている文意表現パターンを個別に無効化する。

【0053】さらに、文書分類装置1は、辞書編集ウィンドウに設けられている新規追加ボタン(図3の右下に示されている新規追加のボタン)が押されると、追加ウィンドウを開き、既存の文意の呼び出しや新規の文意の登録の操作を可能とする。この追加ウィンドウで、例えば、新規ボタンが押されると、その文意タグと対になる文意表現パターンの新規登録を行う。

【0054】例えば、図3では、上記テキスト文書の例において、文意(依頼)に対で登録されていた3パターンを無効化し、追加ウィンドウで既存文意(要望)を呼び出し、追加ウィンドウで呼び出した既存文意(要望)に、先に無効化した文意(依頼)に登録されていた3パターンを登録する処理が示されている。

【0055】文書分類装置1は、ユーザの所望する文意表現パターン辞書の辞書編集が終わると、辞書編集部6が、当該編集結果を文意表現パターン辞書記憶部10の文意表現パターン辞書に保存し、これまでの辞書内容と編集後の辞書内容との辞書の差異を編集履歴として編集履歴記憶部11に保存する。

【0056】辞書編集部6は、ユーザの操作部からの編集履歴参照操作に応じて、編集履歴記憶部11に記憶されている編集履歴を、表示部5に表示し、ユーザが閲覧できるようにする。

【0057】ユーザは、編集履歴を参照して、選択した履歴を、文意表現パターン辞書に反映させて利用することもできる。

【0058】このように編集した辞書を用いて、上記テキスト文書の例において、再度分類を行なうと、文意(依頼)だったデータが、文意(要望)と認識され、分類結果は、文意(推量)が{1, 4, 6, 9, 10, 11}、文意(要望)が{2, 3, 5, 7, 8}という文書の文意別カテゴリ分類が実現される。

【0059】すなわち、図4に示すように、辞書編集部6は、辞書編集指示があると(ステップS101)、辞書編集ウィンドウを起動して、表示部5に表示し(ステ

ップS102)、無効化欄にチェックがあるか判別する(ステップS103)。

【0060】ステップS103で、無効化欄にチェックがないときには、新規追加指示があるか判別し(ステップS104)、新規追加指示がないときには、ステップS101に戻って上記同様に処理する。

【0061】ステップS103で、無効化欄にチェックがあると、辞書編集部6は、無効化欄にチェックの入っている文意の対として登録されている文意表現パターンを個別に無効化する無効化処理を行い(ステップS105)、新規追加指示があるかチェックする(ステップS104)。

【0062】ステップS104で、新規追加指示があると、辞書編集部6は、追加ウィンドウを起動して表示部5に表示させ、既存文意が選択されたかチェックする(ステップS107)。

【0063】ステップS107で、既存文意が選択されると、辞書編集部6は、当該選択された文意を文意表現パターン辞書記憶部10の文意表現パターン辞書から呼び出し(ステップS108)、文意表現パターン辞書に登録する(ステップS109)。

【0064】ステップS107で、既存文意が選択されないときには、辞書編集部6は、新規文意の登録であると判断して、新規文意を文意表現パターン辞書に登録する(ステップS109)。

【0065】文意表現パターン辞書への登録を行うと、辞書編集部6は、保存指示があるかチェックし(ステップS111)、保存指示があると、当該文意表現パターン辞書の編集履歴を編集履歴記憶部11に保存して、処理を終了する(ステップS112)。

【0066】このように、本実施の形態の文書分類装置1は、言語解析部3で、入力された検索要求を含むテキスト文書を、テキスト文書群記憶部2に収集・蓄積された複数のテキスト文書から検索し、言語解析部3で、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、テキスト文書分析部4が、当該カテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書記憶部10の文意表現パターン辞書を参照してテキスト文書を文意別に分析・分類し、当該文意表現パターン辞書を、辞書編集部6で、任意に編集している。

【0067】したがって、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行うことができ、文書分類装置1の利用性を向上させることができる。

【0068】また、本実施の形態の文書分類装置1は、文意表現パターン辞書記憶部10の文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、辞書編集部6が、ユーザの操作に応

じて、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録している。

【0069】したがって、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行うことができ、利用性を向上させることができる。

【0070】さらに、本実施の形態の文書分類装置1は、文意表現パターン辞書記憶部10の文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、辞書編集部6が、ユーザの操作に応じて、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化している。

【0071】したがって、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにすることができ、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性を向上させることができる。

【0072】以上、本発明者によってなされた発明を好適な実施の形態に基づき具体的に説明したが、本発明は上記のものに限定されるものではなく、その要旨を逸脱しない範囲で種々変更可能であることはいうまでもない。

【0073】

【発明の効果】請求項1記載の発明の文書分類装置によれば、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類するに際して、入力された検索要求を含むテキスト文書を、テキスト文書検索手段で、文書蓄積手段に収集・蓄積された複数のテキスト文書から検索し、カテゴリ生成手段で、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、カテゴリ生成手段のカテゴリ分け結果に基づいて、文意分類手段が、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照してテキスト文書を文意別に分析・分類し、当該文意表現パターン辞書を、文意表現パターン辞書編集手段で、任意に編集するので、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行うことができ、利用性を向上させることができる。

【0074】請求項2記載の発明の文書分類装置によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集手段が、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録するので、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行うことができ、利用性を向上させることができる。

【0075】請求項3記載の発明の文書分類装置によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集手段が、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化するので、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにすることができ、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性を向上させることができる。

【0076】請求項4記載の発明の文書分類方法によれば、文書蓄積手段に収集・蓄積された複数のテキスト文書を分析・分類するに際して、テキスト文書検索処理ステップで、入力された検索要求を含むテキスト文書を文書蓄積手段のテキスト文書から検索し、カテゴリ生成処理ステップで、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けして、文意分類処理ステップで、カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書を任意に編集するので、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行うことができ、利用性を向上させることができる。

【0077】請求項5記載の発明の文書分類方法によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書に、意図タグと文意表現パターンを新たに追加登録するので、文意表現パターン辞書に未登録の文意や対応する文意表現パターンを適宜登録して、幅広い文意分類を行うことができ、利用性を向上させることができる。

【0078】請求項6記載の発明の文書分類方法によれば、文意表現パターン辞書に、文書の書き手の意図を示す意図タグと、当該文書の書き手の意図を特定する特徴的な文意表現パターンとを対として登録し、文意表現パターン辞書編集処理ステップで、文意表現パターン辞書に登録されている意図タグと文意表現パターンを個々に無効化するので、文意表現パターンを別文意へ再編成したり、文意表現パターンを削除することなく、無効化とすることで、ユーザが無効化を解除したり、過去履歴を利用する際に再利用できるようにすることができ、ユーザの所望する文意により、より一層適切に文書分類を行って、より一層利用性を向上させることができる。

【0079】請求項7記載の発明の記録媒体によれば、文書蓄積手段に収集・蓄積された複数のテキスト文書を

分析・分類する文書分類方法のプログラム及びデータを記憶する記録媒体に、当該プログラム及びデータとして、入力された検索要求を含むテキスト文書を前記文書蓄積手段のテキスト文書から検索するテキスト文書検索処理ステップと、当該テキスト文書検索結果に基づいて複数のテキスト文書をカテゴリ分けするカテゴリ生成処理ステップと、前記カテゴリ生成処理ステップでのカテゴリ分け結果に基づいて、文意を表現する特徴的な文意表現パターンの登録されている文意表現パターン辞書を参照して前記テキスト文書を文意別に分類する文意分類処理ステップと、前記文意表現パターン辞書を任意に編集する文意表現パターン辞書編集処理ステップと、の各ステップ処理を行う文書分類方法のプログラム及びデータを記録しているので、当該プログラム及びデータをコンピュータ等の情報処理装置に導入することで、分類基準を任意に変更して、ユーザの所望する文意により文書分類を行い、利用性を向上させることのできる文書分類装置及び文書分類方法を実現することができる。

【図面の簡単な説明】

【図1】 本発明の文書分類装置、文書分類方法及び記録

媒体の一実施の形態を適用した文書分類装置の要部ブロック構成図。

【図2】 図1の文書表現パターン辞書記憶部に格納されている文意表現パターン辞書の一例を示す図。

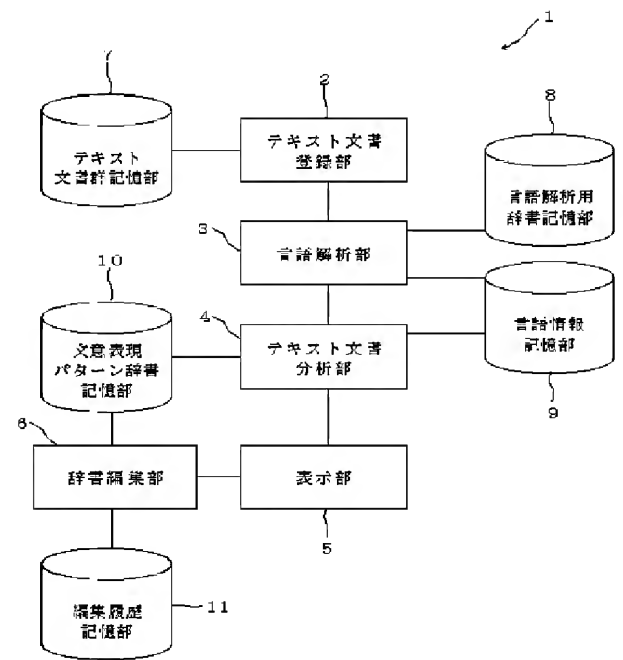
【図3】 図1の表示部に表示される辞書編集ウインドウの一例を示す図。

【図4】 図1の文書分類装置による文意表現パターン辞書編集処理を示すフローチャート。

【符号の説明】

- 1 文書分類装置
- 2 テキスト文書登録部
- 3 言語解析部
- 4 テキスト文書分析部
- 5 表示部
- 6 辞書編集部
- 7 テキスト文書群記憶部
- 8 言語解析用辞書記憶部
- 9 言語情報記憶部
- 10 文意表現パターン辞書記憶部
- 11 編集履歴記憶部

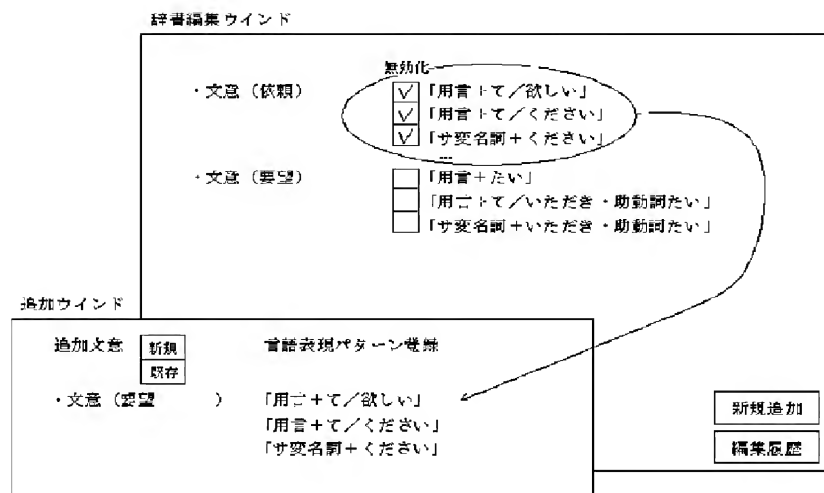
【図1】



【図2】

文意 (種類)	言語表現パターン
・文意 (依頼)	「用言＋て／欲しい」 「用言＋て／ください」 「サ変名詞＋ください」 ...
・文意 (希望)	「用言＋助動詞たい」 「用言＋て／いただき・助動詞たい」 「サ変名詞＋いただき・助動詞たい」 ...
・文意 (否定)	「用言＋助動詞ない」 「サ変名詞＋助動詞＋助動詞ない」 ...
・文意 (疑問)	「用言＋終助詞か」 「文末＝ 記号？」 ...
・文意 (推定)	「用言＋助動詞だろう」 「用言＋副助詞か＋副助詞も／しれない」 「用言＋助動詞そう」 「用言＋助動詞である＋助動詞そう」 ...
...	...

【図3】



【図4】

